



**POLITECNICO**  
**MILANO 1863**

SCUOLA DI INGEGNERIA INDUSTRIALE  
E DELL'INFORMAZIONE

# Systems and Methods for Big and Unstructured Data Project

Author(s): **Gabriele Ginestroni**

**Giacomo Gumiero**

**Lorenzo Iovine**

**Nicola Landini**

**Francesco Leone**

Group Number: **10**

Academic Year: 2022-2023



# Contents

<b>Contents</b>	<b>i</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Problem Specification . . . . .	1
1.2 Assumptions . . . . .	1
<b>2 ER Diagram</b>	<b>3</b>
<b>3 Document Structure</b>	<b>5</b>
3.1 Data Preprocessing . . . . .	5
3.2 Article structure . . . . .	7
3.3 Author structure . . . . .	11
3.4 Attributes Description . . . . .	12
3.4.1 Article . . . . .	12
3.4.2 Author . . . . .	13
<b>4 Commands and Queries</b>	<b>15</b>
4.1 Commands . . . . .	15
4.1.1 Insert a publication in the system . . . . .	15
4.1.2 Insert an author in the system . . . . .	18
4.1.3 Update the number of citations of referenced publications . . . . .	18
4.1.4 Modification of the biography of an author . . . . .	19
4.1.5 Add a publication to its author . . . . .	19
4.2 Queries . . . . .	20
4.2.1 Query 1 . . . . .	20
4.2.2 Query 2 . . . . .	21
4.2.3 Query 3 . . . . .	21
4.2.4 Query 4 . . . . .	22
4.2.5 Query 5 . . . . .	22

4.2.6	Query 6 . . . . .	23
4.2.7	Query 7 . . . . .	24
4.2.8	Query 8 . . . . .	25
4.2.9	Query 9 . . . . .	25
4.2.10	Query 10 . . . . .	26
4.2.11	Query 11 . . . . .	28
4.2.12	Improving Queries Perfomance . . . . .	31

## 5 Conclusions 33

# 1 | Introduction

In this chapter will be presented the problem specification and the hypothesis under which the database is implemented.

## 1.1. Problem Specification

This project aims to build a documental database that handles scientific articles contained in the DBLP bibliography. The focus is on creating a database which allows efficient information retrieval of the articles, including their chapters and images. The main collections analyzed in the project are *Authors* and *Publication* with all their attributes and related objects like: *chapters*, *biographies* and *images*.

## 1.2. Assumptions

1. Articles can be published on a single venue
2. There is no distinction between different types of Publication
3. An author can't work for more than one organization for the same Publication
4. A chapter can contain subsections and these will be considered chapters as well



## 2 | ER Diagram

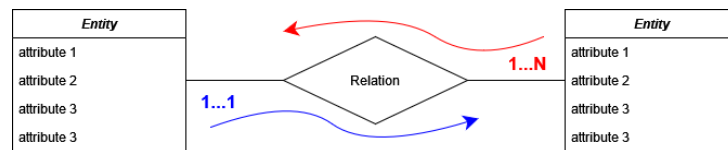


Figure 2.1: ER Diagram Organization

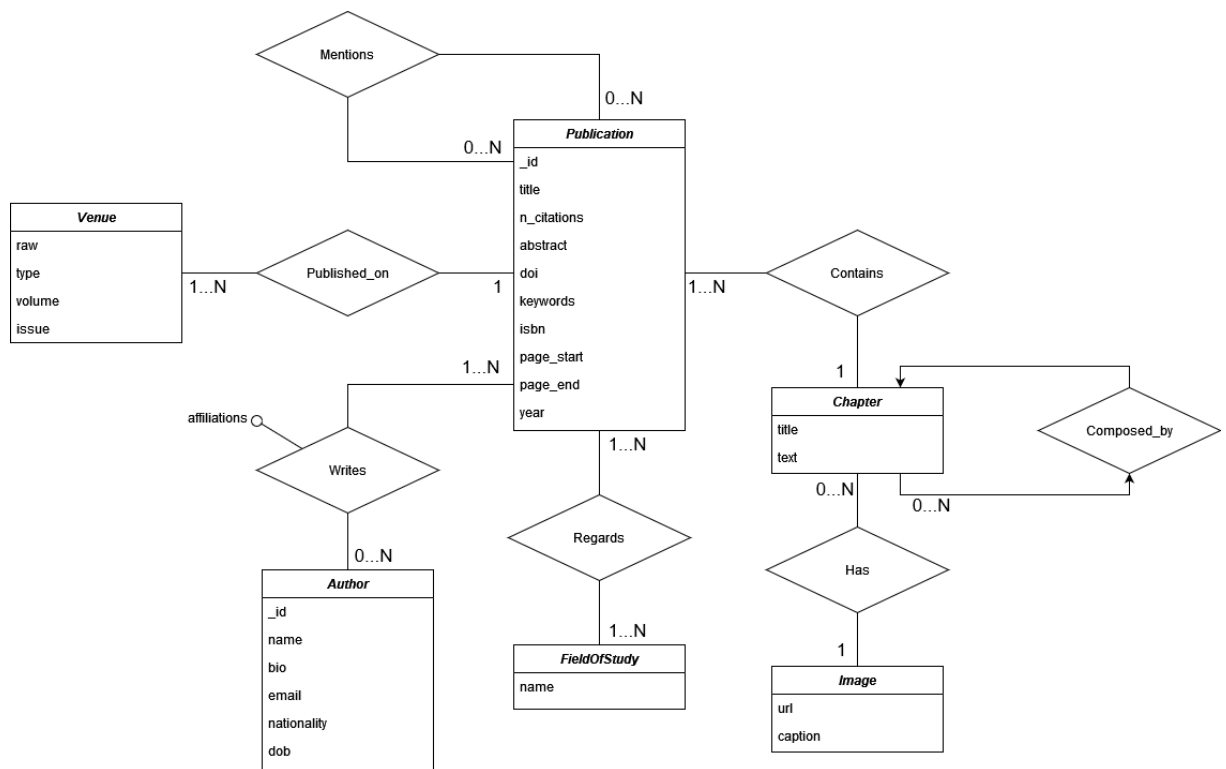


Figure 2.2: ER Diagram

The Entity-Relationship model contains 6 main entities, that are related to each other through various relationships:

- **Publication:** represents all the scientific articles. Its attributes are: *\_id*, *title*, *n\_citations*, *abstract*, *doi*, *keywords*, *isbn*, *page\_start*, *page\_end*, *year* and its organization will be presented later
- **Author:** it is the one who contributed to a publication. Its attributes are: *\_id*, *name*, *bio*, *email*, *nationality*, *dob* (*date of birth*)
- **Venue:** it is where a publication is published or presented. Its attributes are: *raw*, *type*, *volume*, *issue*
- **FieldOfStudy:** this entity represents the subjects of the publication and its attribute is *name*
- **Chapter:** it represents the chapter of a scientific articles. Its attributes are *title*, *text*
- **Image:** it represents the image contained in a specific chapter of the publication. Its attributes are *url*, *caption*

The ER diagram designed contains also the following relationships:

- **Mentions:** is the relationship between a *publication* and another *publication* cited by the first one
- **Published\_on:** is the relationship between a *publication* and its *venue*
- **Writes:** is the relationship between *author* and *publication* which features the affiliation property. We decided to design it with **affiliation** as an attribute of the relationship, due to the fact that it belongs only to a pair of *author* and *publication* and it represents the institute where the author worked for the publication
- **Regards:** is the relationship between a *publication* and its *fields of study*
- **Contains:** is the relationship between *publication* and its *chapters*
- **Composed\_by:** is the relationship between two *chapters*. It represents the relationship created between a chapter and its sections, between a section and its subsection and so on. Note that we used a directed arrow in order to indicate that a section belongs only to a chapter, but a chapter could own more than one section.
- **Has:** is the relationship between *chapter* and its *images*



## 3 | Document Structure

In this section we will describe the structure of our database. We decided to split our dataset in two collections: *authors* and *articles*. The reason behind this choice was to increase the performance and to reduce the spatial complexity. For example if we need to modify the email of an author we just have to change the field in the *author* collection. Whereas, if we just kept one collection, very expensive update query would have been needed (change the modified field in each author's sub-document of the entire article collection).

Furthermore, to avoid redundancies we used manual references to bind the two collections.

In order to show the number of *Authors* and *Articles* in our collections, we performed the two following queries:

```
> db.articles.aggregate([
  { $group: { "_id": 1 , "count": { $sum: 1 }}} ])
< { _id: 1, count: 6381 }
> db.authors.aggregate([
  { $group: { "_id": 1 , "count": { $sum: 1 }}} ])
< { _id: 1, count: 17025 }
```

### 3.1. Data Preprocessing

Starting from the database sample used for the Neo4j part some operations were performed to add the required fields and to exploit better mongoDB functionalities.

Missing *isbn*, *page\_start*, *page\_end* attributes were added with random generations.

Most important, chapters were added from the original paper pdf when available, and parsing was performed by using PHP Library Smalot PDFParser. This library only parses the full text of the papers so the text was then divided into chapters using chapter

titles as a reference to split, thus it was possible to obtain the 2 attributes' title and text for each chapter (for the `text` field, only 10 rows were taken into consideration). A total of 170 pdfs were parsed and that was also used as a pool to pick random chapters for papers that could not be parsed because their pdfs were not available as a free copy on aminers' site.

To add images and captions, a dataset of films and their covers from tmdb was used, for every chapter a random number of images links were added combined with movie titles used as a caption.

The advantage of aminers' dataset was the built-in `_id` attribute for articles and authors and also the existence of an array of references with the list of the referenced `_id` inside every article tuple.

To better exploit this opportunity we believe it was handy to have the 2 collections of articles and authors, so we performed JSON parsing of the dump to generate 2 new JSON files:

- one that contains the articles, with an array of authors that has inside the references to their `_id` objectId
- the other contains the authors with an array of articles that has inside the references to their `_id` objectId.

In the authors' collection the `bio` attribute was randomly generated only when not available, while `email`, `nationality` and `dob` attributes were 100% randomly generated, especially:

- **email:** was generated by concatenating Name.Surname with a pool of most famous mail servers(gmail.com, yahoo.com...)
- **nationality:** was generated by picking randomly from a pool of 20 nation code
- **dob - date of birth:** was generated randomly from 1940 to 1992, so it may be possible that some inconsistencies are present (e.g. author that published before having 20 years old or even before being born)

## 3.2. Article structure

In the following JSON file we can see how the *Article* is structured in our dataset.

**Note:** in some fields we decided not to insert all the content only for ease of read reasons.

**Note:** in `venue`, fields `issue`, `volume` and `publisher` are missing because they are empty in this example.

```
1 {
2   "_id": {
3     "$oid": "53e99f86b7602d9702859fdf"
4   },
5   "title": "Locality Sensitive Outlier Detection: A ranking driven
        approach",
6   "authors": [
7     {
8       "id": {
9         "$oid": "542a4c9fdabfae61d496694e"
10      },
11      "org": "Computer Science and Engineering Department, OhioSU, USA"
12    },
13    {
14      "id": {
15        "$oid": "53f48bc5dabfaea7cd1cce1d"
16      },
17      "org": "Computer Science and Engineering Department, OhioSU, USA"
18    },
19    {
20      "id": {
21        "$oid": "53f44b6fdabfaec09f1dd00d"
22      },
23      "org": "Computer Science and Engineering Department, OhioSU, USA"
24    }
25  ],
26   "n_citation": 60,
27   "abstract": "Outlier detection is fundamental to a variety of database
        and...",
28   "doi": "10.1109/ICDE.2011.5767852",
29   "keywords": [
30     "database point",
31     "ranking scheme",
32     "geometric approach",
33     ...
34  ],
35   "isbn": "978-1-4244-8958-9",
```

```
36 "page_start": "410",
37 "page_end": "421",
38 "year": 2011,
39 "fos": [
40     "Locality-sensitive hashing",
41     "Anomaly detection",
42     "Machine learning",
43     ...
44 ],
45 "venue": {
46     "raw": "ICDE",
47     "type": 0,
48 },
49 "chapters": [
50     {
51         "title": "1. Introduction",
52         "text": "Open Computing Language (OpenCL)[3] is a unified programming
53             ...",
54         "images": [
55             {
56                 "caption": "Giraffada",
57                 "link": "https://image.tmdb.org/t/p/w600_and_h900_bestv2/....jpg"
58             },
59             {
60                 "caption": "A Certain Magical Index: The Miracle Of Endymion",
61                 "link": "https://image.tmdb.org/t/p/w600_and_h900_bestv2/....jpg"
62             },
63             {
64                 "caption": "Ready To Rumble",
65                 "link": "https://image.tmdb.org/t/p/w600_and_h900_bestv2/....jpg"
66             },
67             {
68                 "caption": "Raw Force",
69                 "link": "https://image.tmdb.org/t/p/w600_and_h900_bestv2/....jpg"
70             },
71             {
72                 "caption": "The Departed",
73                 "link": "https://image.tmdb.org/t/p/w600_and_h900_bestv2/....jpg"
74             },
75             {
76                 "caption": "Repeaters",
77                 "link": "https://image.tmdb.org/t/p/w600_and_h900_bestv2/....jpg"
78             }
79         ]
80     }
81 ]
```

```

79  },
80  {
81    "title": "2. The OpenCL Framework",
82    "text": "21 Organization of Our Runtime\nThe target cluster...",
83    "images": [
84      {
85        "caption": "The King Of New York",
86        "link": "https://image.tmdb.org/t/p/w600_and_h900_bestv2/....jpg"
87      },
88      {
89        "caption": "The Ghost Who Walks",
90        "link": "https://image.tmdb.org/t/p/w600_and_h900_bestv2/....jpg"
91      },
92      {
93        "caption": "Fighter In The Wind",
94        "link": "https://image.tmdb.org/t/p/w600_and_h900_bestv2/....jpg"
95      },
96      {
97        "caption": "Wake Up",
98        "link": "https://image.tmdb.org/t/p/w600_and_h900_bestv2/....jpg"
99      },
100     {
101       "caption": "Infierno Blanco",
102       "link": "https://image.tmdb.org/t/p/w600_and_h900_bestv2/....jpg"
103     },
104     {
105       "caption": "Bajo El Mismo Techo",
106       "link": "https://image.tmdb.org/t/p/w600_and_h900_bestv2/....jpg"
107     },
108     {
109       "caption": "Le Dernier Samaritain",
110       "link": "https://image.tmdb.org/t/p/w600_and_h900_bestv2/....jpg"
111     }
112   ]
113 },
114 {
115   "title": "3. Evaluation",
116   "text": "We have implemented the OpenCL runtime and...",
117   "images": [
118     {
119       "caption": "L'hotel Degli Amori Smarriti",
120       "link": "https://image.tmdb.org/t/p/w600_and_h900_bestv2/....jpg"
121     },
122     {

```

```

123     "caption": "Bit",
124     "link": "https://image.tmdb.org/t/p/w600_and_h900_bestv2/....jpg"
125 },
126 {
127     "caption": "Dead Silence",
128     "link": "https://image.tmdb.org/t/p/w600_and_h900_bestv2/....jpg"
129 },
130 {
131     "caption": "Amenazados",
132     "link": "https://image.tmdb.org/t/p/w600_and_h900_bestv2/....jpg"
133 }
134 ]
135 },
136 {
137     "title": "4. Conclusions",
138     "text": "We introduce the design and implementation of...",
139     "images": [
140         {
141             "caption": "Eddie The Eagle",
142             "link": "https://image.tmdb.org/t/p/w600_and_h900_bestv2/....jpg"
143         },
144         {
145             "caption": "Gorenos",
146             "link": "https://image.tmdb.org/t/p/w600_and_h900_bestv2/....jpg"
147         },
148         {
149             "caption": "Pinocchio",
150             "link": "https://image.tmdb.org/t/p/w600_and_h900_bestv2/....jpg"
151         },
152         {
153             "caption": "Le Grand Bazar",
154             "link": "https://image.tmdb.org/t/p/w600_and_h900_bestv2/....jpg"
155         },
156         {
157             "caption": "Gangsters",
158             "link": "https://image.tmdb.org/t/p/w600_and_h900_bestv2/....jpg"
159         },
160         {
161             "caption": "Kodachrome",
162             "link": "https://image.tmdb.org/t/p/w600_and_h900_bestv2/....jpg"
163         }
164     ]
165 }
166 ],

```

```
167 "references": [  
168   {  
169     "$oid": "53e99fddb7602d97028b7e65"  
170   }  
171 ]  
172 }
```

### 3.3. Author structure

In the following JSON file we can see how the *Author* is structured in our dataset.

**Note:** in some fields we decided not to insert all the content only for ease of read reasons.

```
1 {  
2   "_id": {  
3     "$oid": "53f45775dabfaee4dc8162e6"  
4   },  
5   "name": "Guillermo Jorge-Botana",  
6   "articles": [  
7     {  
8       "$oid": "53e99f86b7602d970285a187"  
9     }  
10  ],  
11   "orcid": "0000-0001-5879-6783",  
12   "bio": "Qing-Long Han received the B.Sc. degree in mathematics from the  
13         ...",  
14   "email": "Guillermo.Jorge-Botana@yahoo.com",  
15   "nationality": "de",  
16   "dob": "1974-09-14T00:00:00.000+00:00"  
17 }
```

## 3.4. Attributes Description

In this section we will present all the attributes contained in our filtered dataset.

### 3.4.1. Article

Publication represent the central concept of the system and contains:

- **\_id** is an ObjectId that identifies a publication.
- **title** represents the title of the publication.
- **authors** is an array of subdocuments that contains: ObjectId of the authors of the article and the **org** field which represent the affiliation.
- **n\_citation** is the number of times that the publication has been mentioned.
- **abstract** is a string containing a brief summary of the contents of the paper.
- **doi** Digital Object Identifier is a persistent and standardized identifier.
- **keywords** is an array containing keywords of the publication.
- **isbn** is an identification code of the venue of the publication.
- **page\_start** defines the starting page of the publication.
- **page\_end** defines the last page of the publication.
- **year** represents the year of publication.
- **fos** is an array containing the fields of study of the publication.
- **venue** is a sub-document that represents where a publication is published or presented. This field contains:
  - **raw** is the name or the abbreviation of the venue (regardless the year, issue or volume) in which the publication was presented.
  - **type** indicates the type of the publication.
  - **volume** is the volume of the venue in which the article has been published.
  - **issue** refers to how many times a periodical has been published during that year.
  - **publisher** is the name of the publisher of the article.



- **chapters** is an array of sub-documents containing containing **title** (title of the chapter), **text** (the content), **sections** (that has the same structure of a chapter). They also contains **images** composed by:
  - **caption** the caption of the image.
  - **link** the link to the image.
- **references** set of ObjectIds of the referenced articles.

### 3.4.2. Author

The dataset provides the following author fields:

- **\_id** is an ObjectId that identifies an author.
- **name** is the name of the author.
- **articles** is a set of articles identifier of the publications of the author.
- **orcid** Open Researcher and Contributor ID is a unique identifier for authors of scientific articles.
- **bio** is a string that describes the author.
- **email** is the email address of the author.
- **nationality** is the nationality of the author.
- **dob** is the birth date of the author.



# 4 | Commands and Queries

## 4.1. Commands

We have identified the following **INSERT** and **UPDATE** commands to show the system basic functionalities.

### 4.1.1. Insert a publication in the system

Assuming it is not present in the dataset, we inserted a new document that is a new instance of *Publication*. In order to do that we instantiated 5 different variables: **article\_id** that generates an ObjectId representing the id of the article we're creating; **article\_ref1** and **article\_ref2** that are the ids of two scientific articles cited by this publication; **author1\_id** and **author2\_id** that are the ids of the authors.

**Note:** isbn field is missing

**Note:** type = 1 represents a *Journal*

**Note:** issue field is missing

```

1 article_id = ObjectId()
2 article_ref1 = ObjectId("53e99fe4b7602d97028bf743")
3 article_ref2 = ObjectId("53e99fddb7602d97028bc085")
4 author1_id = ObjectId()
5 author2_id = ObjectId()
6
7 db.articles.insertOne({
8   _id: article_id,
9   title: "An extensive study of C-SMOTE, a Continuous Synthetic Minority
10     Oversampling Technique for Evolving Data Streams",
11   authors:
12     [
13       {id:author1_id, org:"Politecnico di Milano"},
14       {id:author2_id, org:"Politecnico di Milano"}
15     ],
16   n_citation: 3,

```

```
abstract: "Streaming Machine Learning (SML) studies algorithms that
update their models, given an unbounded and often non-stationary flow
of data performing a single pass. Online class imbalance learning is
a branch of SML that combines the challenges of both class imbalance
and concept drift. In this paper, we investigate the binary
classification problem by rebalancing an imbalanced stream of data in
the presence of concept drift, accessing one sample at a time.",
doi: "10.1016/j.eswa.2022.116630",
keywords: ["Evolving Data Stream","Streaming","Concept drift","
Balancing"],
page_start: 39,
page_end: 46,
year: 2022,
fos: ["Computer Science","Stream Reasoning","Big Data"],
venue: {
  raw: "ESA",
  type: 1,
  volume: 196,
  publisher: "Elsevier"
},
chapters:
[
  {
    title:"1. Introduction",
    text:"Nowadays, data abound as a multitude of smart devices, such as
    smartphones, wearables, computers, and Internet of Things (IoT)
    sensors produce massive, continuous, and unbounded flows of data,
    namely data streams. This poses several challenges to Machine
    Learning (ML).",
  },
  {
    title:"2. Background",
    text:"This section is divided into three parts describing the
    different concept drift types and characteristics, the evaluation
    metrics, and the most common approaches to use in class imbalance.",
    subsection: [
      {
        title:"2.1. Concept drift in evolving data streams",
        text:"In this part, we introduce the concept drift phenomenon
        explaining why and how it happens. We explain all its different types
        , forms, and possible speeds of occurrence.",
        subsubsection: [
          {
```

```

45     title:"2.1.1. Concept drift types",
46     text:"Since the generating function is unknown, concept drift is
    unpredictable. In the batch settings, with all the data available,
    it is simple to check and detect if a dataset is not stationary.",
47     images:
48     [
49         {
50             caption:"Fig.1 Representation of the three different types
of concept drift",
51             url:"https://ars.els-cdn.com/content/image/1-s2.0-
S0957417422001208-gr1.jpg"
52         }
53     ]
54     }
55 ]
56 }
57 ]
58 },
59
60 {
61     title:"3. C-SMOTE",
62     text:"This section recalls the description of C-SMOTE, inspired by
the Smote technique, originally presented in Bernardo, Gomes et al.
(2020). C-SMOTE is designed to rebalance an imbalanced data stream,
and it can be pipelined with any SML- model. C-SMOTE stands for
Continuous-Smote, meaning that the new Smote version is applied
continuously.",
63     images:[
64         {
65             caption:"Fig. 2. Architecture of C-SMOTE meta-strategy pipelined
with an Online Learner.",
66             url: "https://ars.els-cdn.com/content/image/1-s2.0-
S0957417422001208-gr5.jpg"
67         }
68     ],
69     subsection:[
70         {
71             title:"3.1. Artificial data streams",
72             text: "To synthetically reproduce the different types of concept
drifts shown in Section 2.1, we choose two of the most commonly used
artificial data generators: SINE1 (Gama et al., 2004) and SEA (Street
& Kim, 2001)."

```

```

75     }
76   ],
77   references: [
78     article_ref1,
79     article_ref2
80   ]
81 ]
82
83 })

```

### 4.1.2. Insert an author in the system

Assuming he is not present in the dataset, we used `insertOne` to create a new instance of *Author*.

**Note:** `author1_id` and `article_id`, refers to the variables instantiated in the previous command (Section: 4.1.1)

```

1 db.authors.insertOne({
2   _id: author1_id,
3   name: "Emanuele Della Valle",
4   orcid: "0000-0002-5176-5885",
5   articles:[
6     article_id
7   ],
8   bio:"Emanuele Della Valle holds a PhD in Computer Science from the
9       Vrije Universiteit Amsterdam and a Master degree in Computer Science
10      and Engineering from Politecnico di Milano. He is associate professor
11      at the Department of Electronics, Information and Bioengineering of
12      the Politecnico di Milano.",
13   email:"emanuele.dellavalle@gmail.com",
14   nationality:"it",
15   dob:ISODate("1975-03-07T00:00:00.000Z")
16 })

```

### 4.1.3. Update the number of citations of referenced publications

With the following snippet of code is possible to increment the `n_citations` field of the *Publications* referenced by the article created in section 4.1.1.

**Note:** in this command we used `updateMany` in order to update both the referenced publications; `update` wasn't enough because only one of the two matching document would have been updated

```
1 db.articles.updateMany(  
2   { $or:[{_id:{$eq:ObjectId("53e99fe4b7602d97028bf743")}}, {_id:{$eq:  
   ObjectId("53e99fddb7602d97028bc085")}}]},  
3   { $inc: { n_citation: 1} }  
4 )
```

#### 4.1.4. Modification of the biography of an author

This command allows to access *Authors* by field **name** in order to append a string to field **bio**

```
1 db.authors.updateMany(  
2   {name: "Emanuele Della Valle"},  
3   [{ $set: { bio: { $concat: [ "$bio", "He recently become full  
   professot at ETH Zurich." ] } } }]  
4 )
```

#### 4.1.5. Add a publication to its author

This command allows to add the new publication (*Publication* creation presented in section 4.1.1), to one of its *Authors*.

**Note:** we assumed **newArticleId** as the identifier of the new publication and that **author1\_id** refers to the variable instantiated in section 4.1.1

```
1 newArticleId = ObjectId()  
2  
3 db.authors.updateOne(  
4   { _id:author1_id },  
5   { $push:{articles:newArticleId}}  
6 )
```

## 4.2. Queries

We have identified the following queries in order to show the system's basic functionalities. For ease of read reasons, we decided to show results obtained using `project` operator, and, in some cases, we showed only some of the results.

### 4.2.1. Query 1

This query returns one publication written after 2013 whose `FieldOfStudy(fos)` contains *'Machine learning'*.

```
1 db.articles.findOne({
2   "$and": [{year: {$gte:2013}}, {"fos":"Machine learning"}]
3 })
```

```
< { _id: ObjectId("53e99f86b7602d970285ab95"),
  title: 'Bootstrapping polarity classifiers with rule-based classification',
  year: 2013,
  fos:
    [ 'Bag-of-words model',
      'Rule-based system',
      'Bootstrapping',
      'Computer science',
      'Artificial intelligence',
      'Classifier (UML)',
      'Margin classifier',
      'Classifier (linguistics)',
      'Linear classifier',
      'Machine learning',
      'Quadratic classifier' ],
  chapters:
    [ { title: '1. Offline Learning of Multiple Regression Model' },
      { title: '2. Analysis on Online Time Complexity' },
      { title: '3. Analysis on Offline Time Complexity' } ] }
```

Figure 4.1: Projection on title, fos, year, chapters.title



### 4.2.2. Query 2

This query returns all the Italian authors born after 1960.

```
1 db.authors.find({
2   $and:[{nationality:'it'},{dob:{$gte:ISODate("1960-01-01")}}]
3 })
```

```
{ _id: ObjectId("53f43107dabfaee02ac8f08b"),
  name: 'WOLFGANG GLÄNZEL',
  articles: [ ObjectId("53e99fc9b7602d97028a40cd") ],
  orcid: '0000-0001-7529-5198',
  bio: 'Hideo Murakami (S\'73-M\'74) received the B.S. degree in electrical engineering from Kanazawa',
  email: 'WOLFGANG.GLÄNZEL@yahoo.com',
  nationality: 'it',
  dob: 1968-10-23T00:00:00.000Z }

{ _id: ObjectId("53f42dbedabfaec09f1197b3"),
  name: 'Murari Mani',
  articles: [ ObjectId("53e99fd0b7602d97028abcd9") ],
  bio: 'Hiromi Narimatsu received the B.Sci. from Tsuda College, Japan, in 2009. She is now in master',
  email: 'Murari.Mani@yandex.ru',
  nationality: 'it',
  dob: 1981-08-15T00:00:00.000Z }
```

### 4.2.3. Query 3

This query returns ten of the articles published by '*Elsevier*' after 2009.

**Description:** this query perform a filtering over the year attribute of the article and on the publisher field of the venue sub-document.

```
1 db.articles.find({
2   $and:[{"venue.publisher":'Elsevier'},{year:{$gte:2009}}]
3 }).limit(10)
```

```
{ _id: ObjectId("53e99fd6b7602d97028b5a7a"),
  title: 'A massively parallel semi-Lagrangian algorithm for solving the transport equation',
  year: 2010,
  fos:
    [ 'Convection-diffusion equation',
      'Monotonic function',
      'Massively parallel',
      'Computer science',
      'Scalar (physics)',
      'Algorithm',
      'Advection',
      'Optical flow',
      'Domain decomposition methods',
      'Scalability' ],
  venue: { raw: 'Procedia Computer Science', publisher: 'Elsevier' },
  chapters:
    [ { title: '1. INTRODUCTION' },
      { title: '2. PROJECT DESCRIPTION' },
      { title: '3. DEMONSTRATION' },
      { title: '4. CONCLUSION AND FUTURE WORK' } ] }
```

Figure 4.2: Projection on title, fos, year, chapters.title, venue.raw, venue.publisher

#### 4.2.4. Query 4

This query returns the top three years sorted by number of publications.

**Description:** the query computes the count by aggregating with respect to the year field. Then groups are sorted by descending order and only the top 3 are kept.

```
1 db.articles.aggregate([
2   {"$group" : {_id:"$year", count:{$sum:1}}},
3   {"$sort": {count:-1}},
4   {"$limit": 3}
5 ]);
```

```
< { _id: 2013, count: 612 }
   { _id: 2010, count: 574 }
   { _id: 2011, count: 561 }
```

#### 4.2.5. Query 5

This query finds all the articles with at least one Stanford affiliation and regarding 'Machine learning' field of study.

**Description:** this query perform a filtering over the fos attribute of the article and on the org field of the authors array of sub-documents.

```
1 db.articles.aggregate([
2   {$match: {"authors.org":{$regex: "Stanford"}}},
3   {$match: {"fos":"Machine learning"}}
4 ])
```

```
< { _id: ObjectId("53e99fe4b7602d97028c5012"),
  title: 'Individualizing generic decision models using assessments as evidence',
  authors:
    [ { id: ObjectId("53f45e83dabfae02ad76294"),
      org: 'Department of Medicine, University of California, San Diego, CA' },
      { id: ObjectId("53f4cb99dabfaeb14f811bc"),
        org: 'Management Science and Engineering Department, Stanford University, Stanford, CA' } ],
  fos:
    [ 'Decision analysis',
      'Decision rule',
      'Data mining',
      'Decision tree',
      'Optimal decision',
      'Computer science',
      'Artificial intelligence',
      'Decision model',
      'Decision field theory',
      'Decision quality',
      'Evidential reasoning approach',
      'Machine learning' ] }
```

Figure 4.3: Projection on authors, fos and title

#### 4.2.6. Query 6

This query returns the top three years sorted by number of distinct authors.

**Description:** the query starts by unwinding the authors array. Then, articles are grouped by year and by author. The distinct count over the year is performed grouping the previous results by year and accumulating 1 for each group. Finally, results are sorted by descending order, keeping only the top 3 years.

```
1 db.articles.aggregate([
2   {$unwind: "$authors"},
3   {$group: {"_id":{"year":"$year", "author":"$authors.id"}}},
4   {$group: {"_id": "$_id.year", "count": { $sum: 1 }}}},
5   {$sort : {"count" : -1}},
6   {$limit : 3 }
7 ]);
```

```
< { _id: 2013, count: 2016 }
    { _id: 2010, count: 1823 }
    { _id: 2011, count: 1750 }
```

#### 4.2.7. Query 7

This query returns the 20 most frequent keywords.

**Description:** the query starts by unwinding the keywords array. The next stage groups with respect to the keywords (case insensitive) and computes the count for each group. Finally, the results are ordered by descending order and limited to the top 20

```
1 db.articles.aggregate([
2   {$unwind: "$keywords" },
3   {
4     $group: {
5       _id: {$toLowerCase: '$keywords'},
6       count: {$sum: 1}
7     }
8   },
9   {$sort : { count : -1}},
10  {$limit : 20}
11 ]);
```

```
< { _id: 'data mining', count: 167 }
    { _id: 'internet', count: 116 }
    { _id: 'computer science', count: 114 }
    { _id: 'information retrieval', count: 100 }
    { _id: 'feature extraction', count: 98 }
    { _id: 'real time', count: 93 }
    { _id: 'algorithm design and analysis', count: 86 }
    { _id: 'neural network', count: 85 }
    { _id: 'indexing terms', count: 83 }
    { _id: 'computational complexity', count: 81 }
    { _id: 'computer architecture', count: 78 }
    { _id: 'hardware', count: 74 }
    { _id: 'mathematical model', count: 73 }
    { _id: 'optimization', count: 73 }
    { _id: 'computational modeling', count: 72 }
    { _id: 'satisfiability', count: 72 }
    { _id: 'machine learning', count: 72 }
    { _id: 'protocols', count: 72 }
    { _id: 'testing', count: 69 }
    { _id: 'application software', count: 67 }
```

### 4.2.8. Query 8

This query, given the title of an article, returns the chapter with the highest number of images.

**Description:** the first stage matches the article given the title. Then an unwind of the chapters array is performed. The project stage is used to compute the variable `imgCount` which stores the number of images contained in each chapter. The results are sorted by descending order and limited to top 1.

```

1 db.articles.aggregate([
2   {"$match" : {title: "Locality Sensitive Outlier Detection: A ranking
   driven approach"}},
3   {"$unwind" : {path: "$chapters"}},
4   {"$project": {
5     "title":1,
6     "chapters":1,
7     "imgCount": { "$size": "$chapters.images" } } },
8   {"$sort": {imgCount:-1}},
9   {"$limit": 1}
10 ]);

```

```

< { _id: ObjectId("53e99f86b7602d9702859fdf"),
  title: 'Locality Sensitive Outlier Detection: A ranking driven approach',
  chapters:
    [ { title: '2. The OpenCL Framework',
      text: '2l Organization of Our Runtime\nThe target cluster architecture consists of a single host node and\nmultiple compute nodes
      images:
        [ { caption: 'The King Of New York',
          link: 'https://image.tmbd.org/t/p/w600_and_h900_bestv2/OxtxltJK6N2aNMb6x4shSdiFwY.jpg' },
          { caption: 'The Ghost Who Walks',
            link: 'https://image.tmbd.org/t/p/w600_and_h900_bestv2/nCJq5mzrzuekS2g0BLc3Q2aYMIQ.jpg' },
          { caption: 'Fighter In The Wind',
            link: 'https://image.tmbd.org/t/p/w600_and_h900_bestv2/hQsnDNn99NDbCvyj9SLxj3hBAkH.jpg' },
          { caption: 'Wake Up',
            link: 'https://image.tmbd.org/t/p/w600_and_h900_bestv2/wv9CJql799ej93urPm5NQScHxrC.jpg' },
          { caption: 'Infierno Blanco',
            link: 'https://image.tmbd.org/t/p/w600_and_h900_bestv2/7IBGsaUXMo12K8i1BTBqxEBuS6y.jpg' },
          { caption: 'Bajo El Mismo Techo',
            link: 'https://image.tmbd.org/t/p/w600_and_h900_bestv2/5SRD0sZ4AnU89L1VuxecD9W2ck2.jpg' },
          { caption: 'Le Dernier Samaritain',
            link: 'https://image.tmbd.org/t/p/w600_and_h900_bestv2/b0wLwVkwGff3088oPyw6fcCzkhH.jpg' } ] },
      imgCount: 7 } ]

```

### 4.2.9. Query 9

This query returns articles citing another article that contains, between its `references`, at least one article written in a Milan University.

**Description:** the first stage loads, by id, the array of cited documents into a new field

called refs. Then a match is used to find in this field all the articles that contain at least 1 author that worked for a Milan university.

**Note:** we used the operator `lookup` to issue a join, in order to access another article instance.

```

1 db.articles.aggregate([
2   {$lookup: {
3     from: "articles",
4     localField: "references",
5     foreignField: "_id",
6     as: "refs"
7   }
8 },
9 {$match: {
10   "refs.authors.org":{$regex: "Milan"}
11 }
12 },
13 {"$project": {
14   "title":1,
15   "refs.title":1,
16   "refs.authors": 1
17 }
18 }
19 ])
```

```

{ _id: ObjectId("53e9a3abb7602d9702cba703"),
  title: 'Adapt cases: extending use cases for adaptive systems',
  refs:
    [ { title: 'Live goals for adaptive service compositions',
        authors:
          [ { id: ObjectId("53f46a7adabfaec09f24fcb5"),
              org: 'Politecnico di Milano, Milano, Italy' },
            { id: ObjectId("53f43983dabfaedce554c415"),
              org: 'Politecnico di Milano, Milano, Italy' } ] ] ] }
```

Figure 4.4: Projection on title, refs.title, refs.authors

#### 4.2.10. Query 10

This query finds the author with the maximum amount of written articles and retrieves his article with the highest number of coauthors.

**Description:** First of all, we compute the `artCount` field which contains the length of

the articles array field. Then we sort in descending order and keep only the first author by artCount. After that, a join is performed to load the articles documents into the new field articles\_doc. A new projection is then used to compute the number of authors of each loaded article. Finally we sort the result in descending order and keep only the top 1 article. The projection is used to display the result in a clearer way.

**Note:** we used the operator `lookup` to issue a join in order to access articles of the author.

```

1 db.authors.aggregate([
2   {"$project": {
3     "_id":1,
4     "name":1,
5     "articles":1,
6     "artCount": { "$size": "$articles" } } },
7   {"$sort": {"artCount": -1}},
8   {"$limit": 1 },
9   {"$lookup": {
10    from: "articles",
11    localField: "articles",
12    foreignField: "_id",
13    as: "articles_doc"}},
14   {"$project": {
15     "name":1,
16     "articles_doc":1,
17     "authCount": { "$size": "$articles_doc.authors" } } },
18   {"$sort": {"authCount": -1}},
19   {"$limit": 1 },
20   {"$project": {
21     "name":1,
22     "articles_doc.title":1,
23     "authCount": 1}}
24 ])
```

```

< { _id: ObjectId("53f48041dabfae963d25910a"),
  name: 'Qiang Yang',
  articles_doc:
    [ { title: 'Adaptive Temporal Radio Maps for Indoor Location Estimation' },
      { title: 'A Whole Page Click Model to Better Interpret Search Engine Click Data.' },
      { title: 'Heterogeneous transfer learning for image clustering via the social web' },
      { title: 'Tracking mobile users in wireless networks via semi-supervised colocalization.' },
      { title: 'Digital Wall: A Power-efficient Solution for Location-based Data Sharing' },
      { title: 'Learning Adaptive Temporal Radio Maps for Signal-Strength-Based Location Estimation' } ],
  authCount: 6 }
```

Figure 4.5: Projection on name, articles\_doc.title, authCount

### 4.2.11. Query 11

This query returns all the articles written by authors whose names combined have all 26 letters of the alphabet.

**Description:** The query is done with the following steps:

- an unwind on articles is performed to obtain an entry for each article written by an author
- results are sorted by authors' name in order to have the list sorted alphabetically for further steps
- for each article, a group is performed and the authors list is pushed into an array
- a string concatenating all the authors names retrieved by article documents, is computed and converted to lowercase
- a map operation is performed to split letter by letter the entries; a filter with a regex to keep only alphabet characters is applied, then an unwind on letters is made
- entries are then grouped by article `_id` and authors list; letters are reduced into an array, then the size of this array is computed and only entries containing all 26 letters are matched
- finally, a join operation is performed on articles to obtain the title

```
1 db.authors.aggregate([
2   {
3     "$unwind" : {path: "$articles"}
4   },
5   { "$sort" : { "name": 1 }},
6   {
7     "$group" : {
8       _id: "$articles",
9       authors: {
10        $push: {
11          $concat: ["$name"]
12        }
13      }
14    },
15  },
16  {
17    $project: {
18      "_id": 1,
19      "authors": "$authors",
20      "fullNames": {
```



```

21     $reduce: {
22         input: "$authors",
23         initialValue: "",
24         in: { $toLowerCase: { $concat : ["$$value", "$$this"]}}
25     }
26 }
27 }
28 },
29 {
30     $project: {
31         "_id": 1,
32         "authors": "$authors",
33         letters: {
34             $filter: {
35                 input: {
36                     $map: {
37                         input: {
38                             $range: [ 0, { "$strLenCP": "$fullNames" } ]
39                         },
40                         in: {
41                             "$substrCP":
42                             [
43                                 "$fullNames",
44                                 "$$this",
45                                 1
46                             ]
47                         }
48                     }
49                 },
50                 cond: {
51                     $regexMatch: {
52                         input: "$$this",
53                         regex: '[a-z]'
54                     }
55                 }
56             }
57         }
58     }
59 },
60 { $unwind: '$letters' },
61 {
62     $group: {
63         _id: {
64             "_id": '$_id',

```

```
65     authors: "$authors"
66   },
67   letters: { $addToSet: '$letters' },
68 }
69 },
70 {
71   $project: {
72     "_id" : 1,
73     "authors": "$authors",
74     "differentLetters": {
75       "$size": "$letters"
76     }
77   }
78 },
79 { $match : {"differentLetters" : 26}},
80 {
81   $lookup: {
82     from: "articles",
83     localField: "_id._id",
84     foreignField: "_id",
85     as: "articles_doc"
86   }
87 },
88 {
89   $project: {
90     "_id": "$_id._id",
91     "title": "$articles_doc.title",
92     "authors": "$_id.authors",
93     "differentLetters": "$differentLetters"
94   }
95 }
96 ])
```

```
{ _id: ObjectId("53e9a042b7602d9702929e98"),
  title: [ 'Building an information retrieval test collection for spontaneous conversational speech' ],
  authors:
    [ 'Bhuvana Ramabhadran',
      'Dagobert Soergel',
      'David S. Doermann',
      'Douglas W. Oard',
      'G. Craig Murray',
      'James Mayfield',
      'Jianqiang Wang',
      'Liliya Kharevych',
      'Martin Franz',
      'Samuel Gustman',
      'Stephanie Strassel',
      'Xiaoli Huang' ],
  differentLetters: 26 }
{ _id: ObjectId("53e9a202b7602d9702b01751"),
  title: [ 'Design principles for developing stream processing applications' ],
  authors:
    [ 'Bugra Gedik',
      'Chitra Venkatramani',
      'Deepak S. Turaga',
      'Henrique Andrade',
      'Jeffrey David Harris',
      'John Cox',
      'Olivier Verscheure',
      'Paul Jones',
      'William Szewczyk' ],
  differentLetters: 26 }
```

Figure 4.6: Projection on `_id`, `title`, `authors`, `differentLetters`

#### 4.2.12. Improving Queries Performance

In our queries, we decided not to use indexes because they weren't necessary to achieve a good performance score. Instead, in bigger and real world databases, could be helpful to create textual indexes on both the collections, in order to perform full-text search on some fields.



## 5 | Conclusions

The documental approach turned out to be very flexible and intuitive, thanks to its affinity with the object-oriented paradigm. Also, this kind of technology allowed us to shape data to match the most frequent operations that could take place in the publication domain. Therefore, queries became very simple and efficient.

Moreover, with our implementation, we tried to reach a good trade-off between performance and spatial complexity, avoiding choices that could lead to critical data duplication problems.